



A distance between populations for one-point crossover in genetic algorithms

Luca Manzoni^a, Leonardo Vanneschi^{a,b}, Giancarlo Mauri^{a,*}

^a Dipartimento di Informatica, Sistemistica e Comunicazione, Università degli Studi di Milano-Bicocca, Milano, Italy

^b ISEGI, Universidade Nova de Lisboa, 1070-312, Lisboa, Portugal

ABSTRACT

Genetic algorithms use transformation operators on the genotypic structures of the individuals to carry out a search. These operators define a neighborhood. To analyze various dynamics of the search process, it is often useful to define a distance in this space. In fact, using an operator-based distance can make the analysis more accurate and reliable than using distances which have no relationship with the genetic operators. In this paper we define a distance which is based on the standard one-point crossover. Given that the population strongly affects the neighborhood induced by the crossover, we first define a crossover-based distance between populations. Successively, we show that it is naturally possible to derive from this function a family of distances between individuals. Finally, we also introduce an algorithm to compute this distance efficiently.

© 2011 Elsevier B.V. All rights reserved.

1. Introduction

Defining a distance based on the neighborhood structure induced by the genetic operators is a basic step for analyzing various dynamics of the search process of Genetic Algorithms (GAs) [1,2]. For instance, it is useful if we want to monitor population diversity (see for instance [2–6]) or if we want to calculate well-known indicators of problem hardness such as fitness–distance correlation (see among others [7,8]). Operator-based distance can make calculating distance and analyzing the search process more accurate [7,9]. Defining a distance, or a function to measure similarity, that is, in some sense “bound” to (or “consistent” with) the genetic operators informally means that if two objects (individuals or populations) are *close* to each other, or similar, one can be transformed into the other with a few applications of the operator(s). This has been recently formalized in [10] and we rely on that definition here.

The aim of this paper is to define a distance that is bound to standard one-point GAs crossover [1,2]. While defining a mutation-based distance can be an easy task for GAs (for instance, Hamming distance is naturally bound to one-point mutation [7]), defining a crossover-based distance can be an issue, mainly because the neighborhood induced by crossover strongly depends on the population where individuals evolve (this difficulty has already been recognized in many references, for instance [7,8]). Nevertheless, defining a crossover-based measure would be extremely important, given that crossover is often the main operator used by GAs to carry out a search.

Here, we solve the problem of the influence of populations on neighborhoods by focusing on the definition of a crossover-based distance between populations. Once that function is defined, we show how it can be used to define a family of distances between individuals.

* Corresponding author.

E-mail addresses: luca.manzoni@disco.unimib.it (L. Manzoni), vanneschi@disco.unimib.it (L. Vanneschi), mauri@disco.unimib.it (G. Mauri).

The definition of a crossover based distance that respects the dynamics of crossover could be used to obtain fitness–distance correlation (fdc) values that are more representative of the real difficulty of the problem. In fact, the fdc value strongly depends on the distance chosen. Hence, a distance that better models the GA dynamics could allow an assessment of the problem difficulty of GA that could be more reliable. Another use of a crossover based distance is the study of the mean distance between a population and all the other possible populations. A population with a low mean distance can be a better initial population for GA because it can allow a better exploration of the search space. Furthermore, the study of a topology over populations can lead to a better understanding of the possible dynamics of GA.

The distance between populations that we introduce is (consistently with the definition given in [10]) strictly related to the minimum number of steps required to transform a population P_1 into another population P_2 by iteratively applying one-point crossover to randomly chosen individuals in P_1 .

Contrary to what one may imagine, we also show that calculating this distance can be computationally cheap and we present an algorithm that performs this calculation in polynomial time with respect to the population size and the number of genes composing the individuals.

The paper is organized as follows. In Section 2 we revise previous and related work. In Section 3, some basic mathematical notions that we use to model GAs are recalled. In Section 4, the model used for computing the proposed distance is studied and some of its general properties are discussed. In Section 5, an alternative and more concise way of representing populations is introduced and an efficient algorithm to compute the proposed distance using this representation is given. Finally, Section 6 discusses and concludes the paper.

2. Previous and related work

The number of contributions published so far aiming at modeling the GAs dynamics is so large that it is impossible to discuss all of them in the restricted space allowed in this paper. The interested reader is referred, for instance, to the numerous papers of Vose and coworkers (for instance [11–17]).

The study of GAs crossover has been carried out in different ways so far. The traditional approach dates back to the early years of the field and it is based on the schema theory [1,2]. Successively, more effective methods for investigating the dynamics induced by crossover have been defined by considering the transition matrix given by it and then studying the Markov process it induces (see for instance [11] and [18]).

A different approach is the one discussed in [19–21], where the topological space induced by crossover is modeled by structures such as hypergraphs and recombination spaces. Although related, the perspective of this paper is different. In fact we aim at defining a distance between populations, which allows us to simplify the structures required for formalizing the model.

The work produced in the last few years by Moraglio and coworkers deserves a particular discussion, given that it is strongly related to the one reported here. In many of its references, among which for instance [22–24], Moraglio gives a geometrical interpretation of many kinds of crossover, including one-point crossover. This allows us to derive distance from operators in a conceptually simple way. One of the many contributions of Moraglio's work stands, in our opinion, in the fact that it sheds light on the importance of studying topologies induced by genetic operators and defining operator-based distances to study evolutionary algorithms (EAs). The approach presented in this paper can be seen, in many senses, as an alternative to the one of Moraglio, the main difference standing in the fact that we focus on the definition of a distance between populations. Also, the mathematical tools we use to model GAs (presented in Section 3) are different from, although related to, the ones used by Moraglio. At present, it is difficult to compare the effectiveness of our work to Moraglio's, given that our work is in its first stage (this paper, in fact, represents the first contribution in the study of crossover-based distances between populations). However, we believe that an alternative approach to existing ones can be interesting for a large part of the EAs community, possibly opening discussions on pros and cons and/or stimulating researchers to investigate possible integrations.

3. Basic notions

In this section some basic notions and some notations that are necessary for the continuation of the paper are introduced.

We denote by $[i, j]$ with $i, j \in \mathbb{N}$ the set $\{i, i + 1, \dots, j - 1, j\} \subseteq \mathbb{N}$. We denote by SC_n the set $\{[i, j] \mid 1 \leq i \leq j \leq n\}$ for a fixed $n \in \mathbb{N}$.

A finite alphabet will be denoted by Σ . The set of all the strings of a given length composed of symbols from Σ is denoted by Σ^n . An element $x \in \Sigma^n$ is denoted by x_1, \dots, x_n . The notation $x_{[i, j]}$ is a shortcut for $x_i, x_{i+1}, \dots, x_{j-1}, x_j$.

Recall that a *lattice* \mathcal{L} is a non-empty set L endorsed with a partial ordering $<_L$ such that for any two elements $a, b \in L$ the *join* $a \vee b$ (i.e., the least upper bound of a and b) and the *meet* $a \wedge b$ (i.e., the greatest lower bound of a and b) operators are uniquely defined in L [25].

A lattice is *bounded* if $\bigvee L$ (i.e., a maximal element for L) and $\bigwedge L$ (i.e., a minimal element for L) exist. A lattice is *complete* if for every subset S of L then both $\bigvee S$ and $\bigwedge S$ exist. Note that every finite lattice is complete.

Given a lattice $\mathcal{L} = (L, <_L)$ a subset O of L is a *lower set* if for all $x \in L$ such that there exists $y \in O$ with $x <_L y$ we have that $x \in O$. The set of all lower sets of a lattice \mathcal{L} is denoted by $\mathcal{O}(\mathcal{L})$ and it is by itself a lattice with respect to set inclusion. See [25] for a reference on lattices.

A Čech closure [26] on a set X is a function $\bar{\cdot} : \mathcal{P}(X) \mapsto \mathcal{P}(X)$ such that:

1. $\bar{\emptyset} = \emptyset$
 2. $\forall A \subseteq X \quad A \subseteq \bar{A}$
 3. $\forall A, B \subseteq X \quad \overline{A \cup B} = \bar{A} \cup \bar{B}$
- (monotonicity)
(additivity)

A Čech closure can be iterated. Define \bar{A}^i with $i \in \mathbb{N}$ as:

$$\bar{A}^i = \begin{cases} \bar{A}^{i-1} & \text{if } i \neq 0 \\ A & \text{otherwise.} \end{cases}$$

When X is finite the function $\llbracket \cdot \rrbracket : \mathcal{P}(X) \mapsto \mathcal{P}(X)$ defined as $\llbracket A \rrbracket = \bigcup_{i \in \mathbb{N}} \bar{A}^i$ is a *Kuratowski closure*. Recall that a *Kuratowski closure* is a Čech closure with idempotency (i.e., $\forall A \subseteq X \llbracket A \rrbracket = \llbracket \llbracket A \rrbracket \rrbracket$) and is one of the ways to define a topology [27].

4. Crossover distance definition

4.1. Crossover relations

In this section we introduce the simplified model for GA with one-point crossover used to define the proposed distance. In this model, populations can be any subset of the set of strings of length n over an alphabet Σ . Hence, we are not considering fixed-size populations and we are not considering the presence of multiple copies of the same individual in the population.

Definition 4.1. A one-point crossover relation R_I is a binary relation over $\Sigma^n \times \Sigma^n$ such that for all $x, y, x', y' \in \Sigma^n$:

$$(x, y)R_I(x', y') \Leftrightarrow \exists k \in [0, n] \text{ s.t. } x' = x_{[1, k]}y_{[k+1, n]} \\ \text{and } y' = y_{[1, k]}x_{[k+1, n]}.$$

In other words, two pairs of elements are in one-point crossover relation when the second pair can be obtained by one-point crossover from the first pair. The relation R_I is reflexive, symmetric but not transitive. It is immediate that its transitive closure is an equivalence relation that partitions $\Sigma^n \times \Sigma^n$ into $\left(\binom{|\Sigma|}{2} + |\Sigma|\right)^n$ equivalence classes (i.e., for every position it is possible to choose a pair of symbols that are not necessarily distinct).

Definition 4.2. A one-point crossover relation R_P over $\mathbf{P} = \mathcal{P}(\Sigma^n)$ is a relation such that for all $P_1, P_2 \in \mathbf{P}$:

$$P_1 R_P P_2 \Leftrightarrow \forall x' \in P_2 \exists y' \in \Sigma^n \exists x, y \in P_1 \\ \text{s.t. } (x, y)R_I(x', y').$$

In other words, two subsets of Σ^n are in relation if and only if every element of the second subset can be obtained by crossover from elements of the first subset. Notice that there are some assumptions in this model of crossover. First of all not all elements of the first population need to contribute to obtain the second population. Furthermore, only one of the offspring need to be inserted into the resulting population. Finally, the current model does not take into account the size of the population that can be obtained (e.g., $\Sigma^n R_P \emptyset$).

Example 4.1. Consider the following two populations:

$$P_1 = \{0110, 0101, 0011\} \\ P_2 = \{0111, 0001, 0000\}.$$

We have that $P_2 R_P P_1$ since all of the elements of P_1 can be obtained with one application of one-point crossover to a pair of elements of P_2 . In fact, 0110 can be obtained from 0111 and 0000, 0101 can be obtained from 0111 and 0001, while 0011 can be obtained from 0000 and 0111. It is not true that $P_1 R_P P_2$ since 0000 cannot be obtained from one application of one-point crossover from any pair of elements of P_1 .

The relation R_P is reflexive. The Example 4.1 shows that R_P is not a symmetric relation. Neither is the relation transitive. Note that, by definition, if $P_1 R_P P_2$ then for all $P'_2 \subseteq P_2$ and for all $P'_1 \supseteq P_1$ it is true that $P'_1 R_P P'_2$.

The main idea that will be carried on is to define a Čech closure $\bar{\cdot}$ such that for any population P we have that \bar{P}^i is the set of populations that can be obtained after i generations using only crossover as a genetic operator. In this way it is possible to define a closure $\llbracket \cdot \rrbracket$ such that for any two populations P_1 and P_2 :

- P_2 can be obtained using only crossover from P_1 iff $P_2 \in \llbracket \{P_1\} \rrbracket$.
- The minimal $k \in \mathbb{N}$ such that $P_2 \in \bigcup_{i=0}^k \bar{P}_1^i$ is the minimum number of generations needed to obtain P_2 from P_1 .

Such a closure can be defined in the following way:

Definition 4.3. The *crossover closure* is a function $\bar{\cdot} : \mathcal{P}(\mathbf{P}) \mapsto \mathcal{P}(\mathbf{P})$ defined, for every $A \subseteq \mathbf{P}$, as :

1. When $A = \emptyset$, $\bar{A} = \emptyset$.
2. When $A = \{P\}$, $\bar{P} = \{P' \in \mathbf{P} \mid P R_P P'\}$.
3. Otherwise, $\{P_1, P_2, \dots, P_k\} = \bigcup_{i=1}^k \bar{P}_i$.

The first property we need to prove is that $\overline{\cdot}$ is a Čech closure. In this way we will be able to use the properties of Čech closures when needed.

Proposition 4.1. *The crossover closure is a Čech closure.*

Proof. The closure of \emptyset is \emptyset by definition. Since the relation R_p is reflexive, for all $A \subseteq \mathbf{P}$ it is immediate that $A \subseteq \overline{A}$. The property of additivity holds by the definition of crossover closure. \square

Furthermore, for all $P_1, P_2 \in \mathbf{P}, P_2 \in \overline{P_1}^k$ iff there exist

$$P_1 = Q_0, Q_1, \dots, Q_{k-1}, Q_k = P_2 \in \mathbf{P}$$

such that $Q_i R_p Q_{i+1}$ for all $i \in [0, k-1]$. In other words, we are requiring that one application the closure $\overline{\cdot}$ effectively represents one generation.

Proposition 4.2. *For all $P_1, P_2 \in \mathbf{P}$ and for all $k \in \mathbb{N}, P_2 \in \overline{P_1}^k$ (Property 1) iff there exists $Q_0, \dots, Q_k \in \mathbf{P}$ with $Q_0 = P_1, Q_k = P_2$ and such that $Q_i R_p Q_{i+1}$ for $i \in [0, k-1]$ (Property 2).*

Proof. Suppose that Property 1 holds. We prove by induction on k that Property 2 follows. When $k = 0$ it is immediate that for any $P_1, P_2 \in \mathbf{P} P_2 \in \overline{P_1}$ iff $P_1 = P_2$ which immediately gives $P_1 R_p P_1$ by the reflexivity of R_p .

Suppose that Property 1 implies Property 2 up to k . We prove that the implication is also true for $k+1$. Take any $P_1, P_2 \in \mathbf{P}$ such that $P_2 \in \overline{P_1}^{k+1}$. There exists a population $P' \in \overline{P_1}^k$ such that $P' R_p P_2$ by definition of $\overline{\cdot}$. By induction hypothesis there exists $P_1 = Q_0, Q_1, \dots, Q_k = P'$ such that $Q_i R_p Q_{i+1}$ for all $i \in [0, k-1]$. The sequence $P_1 = Q_0, \dots, Q_k = P', Q_{k+1} = P_2$ is such that $Q_i R_p Q_{i+1}$ for all $i \in [0, k]$ proving that Property 2 holds for $k+1$.

Now assume Property 2. We prove that it implies Property 1. When $k = 0$ the statement is vacuously true.

Suppose that Property 2 implies Property 1 up to k . We prove that the implication also holds for $k+1$. Take any $P_1, P_2 \in \mathbf{P}$ such that there exists $P_1 = Q_0, Q_1, \dots, Q_k, Q_{k+1} = P_2$ such that Property 2 holds. By induction hypothesis $Q_k \in \overline{P_1}^k$. By the definition of $\overline{\cdot}$ we have that $P_2 \in \overline{P_1}^{k+1}$. Thus, Property 1 holds for $k+1$. \square

Proposition 4.2 shows that iterating the closure of a set of populations k times is equivalent to collecting all the populations reachable from the considered set in at most k generations using one-point crossover.

4.2. The structure of the closure

We study the structure of \overline{P} for all $P \in \mathbf{P}$. In many cases a chain of sets of populations $\{P\} \subseteq \overline{P} \subseteq \overline{P}^2 \subseteq \dots$ could be substituted by a chain of populations $P \subseteq P' \subseteq P'' \subseteq \dots$ that is more easily tractable.

Definition 4.4. Let $\mu_P : \mathcal{P}(\mathbf{P}) \mapsto \mathbf{P}$ be defined as:

$$\mu_P(\{P_1, \dots, P_k\}) = \bigcup_{i=1}^k P_i.$$

Proposition 4.3. *For all $P \in \mathbf{P}, \overline{P}$ is the lattice $(\mathcal{P}(\mu_P(\overline{P})), \subseteq)$.*

Proof. It is immediate that for all $P' \in \overline{P} \implies P' \subseteq \mu_P(\overline{P})$. We only need to prove the inverse implication. Note that $\{x\} \in \overline{P} \iff \{x\} \subseteq \mu_P(\overline{P})$ and that $\emptyset \in \overline{P}$. Then we only need to prove that \overline{P} is closed under finite union. Let $P_1, P_2 \in \overline{P}$ then for all $x' \in P_1 \cup P_2$ there exists $y' \in \Sigma^n$ and $x, y \in P$ such that $(x, y) R_l(x', y')$ since either $x' \in P_1$ or $x' \in P_2$. This means that $P_1 \cup P_2 \in \overline{P}$.

Since \overline{P} is the power set of $\mu_P(\overline{P})$ it is a lattice with respect to the set inclusion ordering with union and intersection as join and meet operations. \square

Note that \overline{P} has $\mu_P(\overline{P})$ as a maximal element and \emptyset as a minimal element.

Proposition 4.4. *For all $P \in \mathbf{P}$ and for all $i \in \mathbb{N}, \overline{P}^i$ is a lattice.*

Proof. $\{P\}$ is a lattice. \overline{P} is a lattice by Proposition 4.3. Suppose that \overline{P}^i is a lattice. We prove that \overline{P}^{i+1} is also a lattice. Due to the additivity property of Čech closures we have that $A \subseteq B \implies \overline{A} \subseteq \overline{B}$. Also, since for all $P' \in \overline{P}^i, P' \subseteq \mu_P(\overline{P}^i)$ we have that

$$\overline{P}^{i+1} = \bigcup_{P' \in \overline{P}^i} \overline{P'} \subseteq \overline{\mu_P(\overline{P}^i)}.$$

The other direction of the inclusion is immediate because $\mu_P(\overline{P}^i) \in \overline{P}^i$. Since $\mu_P(\overline{P}^i) \in \mathbf{P}$ we have that the closure of $\mu_P(\overline{P}^i)$ is a lattice. \square

As a direct corollary given by the proof of the previous proposition, we have that:

Corollary 4.5. For all $P \in \mathbf{P}$ and for all $i \in \mathbb{N}$ with $i > 0$, $\overline{\{P\}}^i$ is the lattice of $\mathcal{P}(\mu_P(\overline{\{P\}}^i))$ ordered by set inclusion.

To a sequence $\{P\} \subseteq \overline{\{P\}} \subseteq \overline{\{P\}}^2 \subseteq \dots$ we can associate the sequence $\mu_P(\{P\}) \subseteq \mu_P(\overline{\{P\}}) \subseteq \mu_P(\overline{\{P\}}^2) \subseteq \dots$ that is also equivalent to the sequence $S_0(P) \subseteq S_1(P) \subseteq S_2(P) \subseteq \dots$ where $S_0(P) = P$ and $S_i(P) = \mu_P(\overline{\{S_{i-1}(P)\}})$ when $i > 0$. This means that, in order to know if a population P' exists inside $\overline{\{P\}}^i$, all we have to do is to determine if P' is a subset of $S_i(P)$. Also, to know if there exists a population P' inside $\overline{\{P\}}^i$ such that a certain element x is in P' , all we have to do is seeking if $x \in S_i$. Now, we study the properties that hold in both representations.

Definition 4.5. Let $P \in \mathbf{P}$. Then for all $i \in \mathbb{N}$ the set $S_i(P) \in \mathbf{P}$ is defined as:

$$S_i(P) = \begin{cases} P & \text{if } i = 0 \\ \mu_P(\overline{\{S_{i-1}(P)\}}) & \text{otherwise.} \end{cases}$$

The function $next : \mathbf{P} \mapsto \mathbf{P}$ is defined as $next(P) = S_1(P)$.

The first property that can be reported from a representation to the other is the presence of a population in a closure. In this case the sentence $P_2 \in \overline{\{P_1\}}^i$ simply becomes $P_2 \subseteq S_i(P_1)$. First of all we need the following proposition linking $S_i(P)$ with the Čech closure.

Proposition 4.6. For all $P \in \mathbf{P}$ and for all $i \in \mathbb{N}$ with $i > 1$, $\overline{\{P\}}^i = \overline{\{S_{i-1}(P)\}}$.

Proof. By induction on i , consider $i = 1$, then $\overline{\{P\}} = \overline{\{S_0(P)\}}$ by the definition of $S_0(P)$. Consider the equivalence true for i , we are going to prove it for $i + 1$. $\overline{\{P\}}^{i+1} = \overline{\{P\}}^i = \overline{\overline{\{S_{i-1}(P)\}}}$. By the proof of Proposition 4.4 we have that $\overline{\overline{\{S_{i-1}(P)\}}} = \mu_P(\overline{\{S_{i-1}(P)\}})$ that, by definition of $S_i(P)$, is equal to $\overline{S_i(P)}$. \square

The desired corollary is then the following one.

Corollary 4.7. For all $P_1, P_2 \in \mathbf{P}$ and for all $i \in \mathbb{N}$ with $i > 0$, $P_2 \in \overline{\{P_1\}}^i$ iff $P_2 \subseteq S_i(P_1)$.

Proof. $P_2 \in \overline{\{P_1\}}^i$ is equivalent to $P_2 \in \overline{\{S_{i-1}(P_1)\}}$. Since for $i > 0$, $\overline{\{P\}}^i$ contains all the subsets of $\mu_P(\overline{\{P\}}^i)$, $P_2 \in \overline{\{S_{i-1}(P_1)\}}$ is equivalent to $P_2 \subseteq \mu_P(\overline{\{S_{i-1}(P_1)\}}) = S_i(P_1)$. \square

Proposition 4.8. Let $P_1, P_2 \in \mathbf{P}$ such that there exists $i \in \mathbb{N}$ with $P_2 \subseteq S_i(P_1)$. Then the following holds:

$$\min\{i \in \mathbb{N} \mid P_2 \subseteq S_i(P_1)\} = \max\{\min\{i \in \mathbb{N} \mid \{x\} \subseteq S_i(P_1)\} \mid x \in P_2\}.$$

Proof. Let the two sides of the equation be called ℓ_1 and ℓ_2 respectively and suppose ℓ_1 and ℓ_2 are different from 0 (the proposition is immediately proved otherwise). Suppose $\ell_1 < \ell_2$. This is impossible since when $P_2 \subseteq S_{\ell_1}(P_1)$ we have also that $\{x\} \subseteq S_{\ell_1}(P_1)$ for all $x \in P_2$. Hence, $\ell_1 \geq \ell_2$. Suppose $\ell_1 > \ell_2$. This is also impossible since we have that $\{x\} \subseteq S_{\ell_2}(P_1)$ for all $x \in P_2$. By Corollary 4.7 this means that $\{x\} \in \overline{\{P_1\}}^{\ell_2}$ for all $x \in P_2$. By the lattice structure of the closure of $\overline{\{P_1\}}^{\ell_2}$ we also have that $\bigcup_{x \in P_2} \{x\} = P_2 \in \overline{\{P_1\}}^{\ell_2}$. This means that $P_2 \in S_{\ell_2}(P_1)$. Hence $\ell_1 = \ell_2$. \square

Now we prove that computing these minimal values on the sets $\overline{\{P\}}^i$ is similar to computing them on the sets $S_i(P)$, hence we can use the latter sets to obtain information on the former ones.

Proposition 4.9. For all $P_1, P_2 \in \mathbf{P}$ such that there exists $i \in \mathbb{N}$ with $P_2 \in \overline{\{P_1\}}^i$ the following holds:

$$\min\{i \in \mathbb{N} \mid P_2 \in \overline{\{P_1\}}^i\} = \begin{cases} 1 & \text{if } P_2 \subset P_1 \\ 0 & \text{if } P_1 = P_2 \\ \min\{i \in \mathbb{N} \mid P_2 \subseteq S_i(P_1)\} & \text{otherwise.} \end{cases}$$

Proof. In the first case we have that $P_2 \subset P_1$ implies that the value of i must be at least 1. It is 1 because of the fact that the closure $\overline{\{P_1\}}$ is a lattice of subsets that contains P_1 and, consequently, P_2 . The second case is immediate since $P_1 \in \{P_1\}$. The third case is provided by Corollary 4.7 if we assume that the minimal i is not 0. Since $P_2 \subseteq S_0(P_1)$ is equivalent to $P_2 \subseteq P_1$ we have that this conditions are covered by the other two cases. \square

4.3. Distance definition

First of all, we define a quasi-metric that will be successively used to define a metric over \mathbf{P} .

Recall that a quasi-metric is a function d such that:

1. For all x, y , $d(x, y) \geq 0$ and $d(x, y) = 0 \Leftrightarrow x = y$.
2. For all x, y, z , $d(x, y) \leq d(x, z) + d(y, z)$.

Note that iterating the closure operator $\bar{\cdot}$, we always reach a fixed point after a finite number of steps (i.e., there exists $k \in \mathbb{N}$ such that $\bar{\cdot}^k$ and $\bar{\cdot}^{k+1}$ are the same function), since the Čech closure is monotone and the domain \mathbf{P} is finite. In fact, the fixed point is necessarily reached after at most $|\mathbf{P}|$ iterations.

Let k^* be the integer $\min\{k \in \mathbb{N} \mid \forall U \subseteq \mathbf{P} : \bar{U}^k = \bar{U}^{k+1}\}$ (i.e., the first $k \in \mathbb{N}$ that allows any possible iteration of $\bar{\cdot}$ to reach a fixed point). Note that for any two populations P_1, P_2 if $P_2 \notin \overline{\{P_1\}}^{k^*}$ then P_2 is not reachable by crossover from P_1 .

Definition 4.6. Let $f_P : \mathbf{P} \times \mathbf{P} \mapsto \mathbb{R}_+$ be defined as:

$$f_P(P_1, P_2) = \begin{cases} \min\{k \in \mathbb{N} \mid P_2 \in \overline{\{P_1\}}^k\} & \text{if } P_2 \in \overline{\{P_1\}}^{k^*} \\ k^* + 1 & \text{otherwise.} \end{cases}$$

Proposition 4.10. The function f_P is a quasi-metric.

Proof. It is immediate that f_P is always not negative. Also, $f_P(P_1, P_2) = 0$ iff $P_2 \in \{P_1\}$ (i.e., iff $P_1 = P_2$).

For the sake of argument, suppose that the triangle inequality does not hold. Then there exist $P_1, P_2, P' \in \mathbf{P}$ such that $f_P(P_1, P') + f_P(P', P_2) < f_P(P_1, P_2)$. Without loss of generality suppose all considered values of f_P less than $k^* + 1$. By Proposition 4.2 there exist $P_1 = Q_0, \dots, Q_{f_P(P_1, P')} = P'$ and $P' = S_0, \dots, S_{f_P(P', P_2)} = P_2$ such that $Q_i R_P Q_{i+1}$ for all $i \in [0, f_P(P_1, P') - 1]$ and $S_i R_P S_{i+1}$ for all $i \in [0, f_P(P', P_2) - 1]$. It is possible to concatenate the two sequences to obtain $P_1 = Q_0, \dots, Q_{f_P(P_1, P')} = P' = S_0, \dots, S_{f_P(P', P_2)} = P_2$. Also by Proposition 4.2 we have that $P_2 \in \overline{\{P_1\}}^k$ with $k = f_P(P_1, P') + f_P(P', P_2)$. By the definition of f_P we have that $f_P(P_1, P_2) \leq f_P(P_1, P') + f_P(P', P_2)$, contradicting the initial assumption of triangle inequality to be false. \square

From a quasi-metric it is immediate to define a metric by summing to f_P itself with swapped arguments in order to obtain symmetry.

Definition 4.7. Let $d_P : \mathbf{P} \times \mathbf{P} \mapsto \mathbb{R}_+$ be defined as:

$$d_P(P_1, P_2) = \frac{1}{2} (f_P(P_1, P_2) + f_P(P_2, P_1)).$$

Note that for every population it is possible to define a distance between individuals.

Definition 4.8. Let $P \in \mathbf{P}$. Then the function $d_P^P : \Sigma^n \times \Sigma^n \mapsto \mathbb{R}_+$ is defined as:

$$d_P^P(x, y) = d_P((P \setminus \{x\}) \cup \{y\}, (P \setminus \{y\}) \cup \{x\}).$$

Proposition 4.11. For all $P \in \mathbf{P}$, the function d_P^P is a distance.

Proof. Both symmetry and the triangle inequality are inherited from the fact that d_P is a distance. The only property that need proving is that for all $x, y \in \Sigma^n$, $x = y \Leftrightarrow d_P^P(x, y) = 0$. This is immediate since $(P \setminus \{x\}) \cup \{y\}$ can be equal to $(P \setminus \{y\}) \cup \{x\}$ only when $x = y$. \square

Note that all the steps from Definition 4.3 are not dependent on the explicit definition of the crossover relation. In fact, all the definitions and propositions remain valid also for any other relation. Their extension to other kinds of crossover is then immediate.

Also, note that it is possible to use Corollary 4.9 and Proposition 4.8 to decompose the computation of the distance between populations into a series of computations of $\min\{i \in \mathbb{N} \mid \{x\} \in S_i(P)\}$ for some individual x and population P . Therefore an efficient method to carry on this computation translates immediately in an efficient way of computing the proposed distance.

5. A concise model for populations

In this section we define a succinct representation for populations.

Definition 5.1. We define as *crossover granules* (SC_n, \subseteq) the set $SC_n = \{[i, j] \mid 1 \leq i \leq j \leq n\}$ ordered by set inclusion.

Proposition 5.1. (SC_n, \subseteq) is a lattice.

Proof. Let $[i, j], [h, k] \in \text{SC}_n$. We show that both $[i, j] \wedge [h, k]$ and $[i, j] \vee [h, k]$ exist and are $[\max\{i, h\}, \min\{j, k\}]$ (that will be denoted by $[M_1, m_1]$) and $[\min\{i, h\}, \max\{j, k\}]$ (that will be denoted by $[m_2, M_2]$) respectively.

It is immediate that $[M_1, m_1] \subseteq [i, j]$ and $[M_1, m_1] \subseteq [h, k]$. It is necessary to prove that every other $[\ell_1, \ell_2] \in \text{SC}_n$ such that $[\ell_1, \ell_2] \subseteq [i, j]$ and $[\ell_1, \ell_2] \subseteq [h, k]$ is also such that $[\ell_1, \ell_2] \subseteq [M_1, m_1]$. Suppose $[M_1, m_1] \neq \emptyset$ otherwise the property is vacuously true. For the sake of argument suppose that there exists $a \in [\ell_1, \ell_2]$ such that $a \notin [M_1, m_1]$. This means that either $a < M_1$ or $a > m_1$. Without loss of generality suppose we are in the first case. Then either $a < i$ or $a < h$. This means that $a \notin [i, j]$ or $a \notin [h, k]$ in contradiction with one of the hypotheses. Thus, $[i, j] \wedge [h, k] = [M_1, m_1]$.

It is also immediate that $[m_2, M_2] \supseteq [i, j]$ and $[m_2, M_2] \supseteq [h, k]$. It is necessary to prove that every other $[\ell_1, \ell_2] \in \text{SC}_n$ such that $[\ell_1, \ell_2] \supseteq [i, j]$ and $[\ell_1, \ell_2] \supseteq [h, k]$ is also such that $[\ell_1, \ell_2] \supseteq [m_2, M_2]$. For the sake of argument suppose that there exists $a \in [m_2, M_2]$ such that $a \notin [\ell_1, \ell_2]$. By definition $m_2 \leq a \leq M_2$. Since $a \notin [\ell_1, \ell_2]$ then either $\ell_1 > a \geq m_1$ or $\ell_2 < a \leq M_2$. Without loss of generality suppose that we are in the first case. Then $m_2 \notin [\ell_1, \ell_2]$ but $m_2 = i$ or $m_2 = h$. This means that $[\ell_1, \ell_2] \not\supseteq [i, j]$ or $[\ell_1, \ell_2] \not\supseteq [h, k]$, negating one of the hypotheses. Thus $[i, j] \vee [h, k] = [m_2, M_2]$.

Since both the meet and the join exist for every and are unique for every pair of elements, (SC_n, \subseteq) is a lattice. \square

Definition 5.2. Let $x \in \Sigma^n$, $[i, j] \in \text{SC}_n$ and $P \in \mathbf{P}$. We say that $x(i, j)$ is represented in P iff there exists $y \in P$ such that $y_{[i,j]} = x_{[i,j]}$.

Note that if $x(i, j)$ is represented in P it is also true that for all $h \geq i$ and for all $k \leq j$, $x(h, k)$ is represented in P . Also note that if $x(1, n)$ is represented in P then $x \in P$.

It is now possible to define the concept of representation for a population.

Definition 5.3. Fix $x \in \Sigma^n$. We define $r_x : \mathbf{P} \mapsto \mathcal{P}(\text{SC}_n)$ as $r_x(P) = \{[i, j] \in \text{SC}_n \mid x(i, j) \text{ is represented in } P\}$.

We now prove that $r_x(P)$ is always a lower set of SC_n .

Proposition 5.2. For all $P \in \mathbf{P}$ and for all $x \in \Sigma^n$, $r_x(P) \in \mathcal{O}(\text{SC}_n)$.

Proof. Consider $[i, j] \in r_x(P)$. By definition $x(i, j)$ is represented in P . This means that also all $x(h, k)$ with $h \geq i$ and $k \leq j$ are represented in P . In other words, $[h, k] \in r_x(P)$. Recall that the elements of SC_n in the form $[h, k]$ with $h \geq i$ and $k \leq j$ are all the elements of SC_n with $[h, k] \subseteq [i, j]$. Since they are in $r_x(P)$ we have that it is a lower set. \square

Now we define a function on $\mathcal{O}(\text{SC}_n)$ that can be used to “mimic” the Čech closure defined over \mathbf{P} .

Definition 5.4. Let $U \in \mathcal{O}(\text{SC}_n)$. We define $\mu_{\text{SC}} : \mathcal{O}(\text{SC}_n) \mapsto \mathcal{O}(\text{SC}_n)$ as

$$\begin{aligned} \mu_{\text{SC}}(U) &= \{[i, j] \in \text{SC}_n \mid \exists [h_1, k_1], [h_2, k_2] \in U \text{ s.t.} \\ &\quad [i, j] = [h_1, k_1] \vee [h_2, k_2] = [h_1, k_1] \cup [h_2, k_2]\}. \end{aligned}$$

Note that $\mu_{\text{SC}}(U)$ can also be formulated as:

$$\bigcup_{\substack{[h_1, k_1], [h_2, k_2] \in U \\ k_1 \geq h_2 - 1 \text{ or} \\ k_2 \geq h_1 - 1}} [h_1, k_1] \vee [h_2, k_2].$$

We are now going to state and prove the main result that allow us to work on the lower set of the lattice SC_n instead of \mathbf{P} .

Theorem 5.3. For all $x \in \Sigma^n$, the following diagram commutes:

$$\begin{array}{ccc} \mathbf{P} & \xrightarrow{\text{next}} & \mathbf{P} \\ \downarrow r_x & & \downarrow r_x \\ \mathcal{O}(\text{SC}_n) & \xrightarrow{\mu_{\text{SC}}} & \mathcal{O}(\text{SC}_n) \end{array}$$

In other words, $r_x \circ \text{next} = \mu_{\text{SC}} \circ r_x$.

Proof. Fix $P \in \mathbf{P}$ and $x \in \Sigma^n$. Consider $z, v \in P$. Also, recall that $r_x(P)$ can be seen as $\bigcup_{y \in P} r_x(\{y\})$. Firstly, we are going to prove that $r_x(\text{next}(P)) \subseteq \mu_{\text{SC}}(r_x(P))$. Let $y \in \Sigma^n$ be such that there exists $w \in \Sigma^n$ such that $(z, v)R_I(y, w)$. Then $y \in \text{next}(P)$. Consider $r_x(\{y\})$. Let $[i, j] \in r_x(\{y\})$. This means that $y_{[i,j]} = x_{[i,j]}$ there can be two cases:

1. Either $[i, j] \in r_x(\{z\})$ or $[i, j] \in r_x(\{v\})$. In this case $[i, j] \in \mu_{\text{SC}}(r_x(\{z\}))$ or $[i, j] \in \mu_{\text{SC}}(r_x(\{v\}))$ (since μ_{SC} is monotone).
2. Neither $[i, j] \in r_x(\{z\})$ nor $[i, j] \in r_x(\{v\})$. In this case, by definition of crossover there we must have k such that $i \leq k < j$, $z_{[i,k]} = x_{[i,k]}$ and $v_{[k+1,j]} = x_{[k+1,j]}$ (or the same with z and v swapped). Hence $[i, k] \in r_x(\{z\})$ and $[k+1, j] \in r_x(\{v\})$. We have that $[i, k] \vee [k+1, j] = [i, k] \cup [k+1, j] = [i, j]$. Therefore $[i, j] \in \mu_{\text{SC}}(r_x(\{z\}) \cup r_x(\{v\}))$.

Combining both cases for all $[i, j] \in r_x(\text{next}(P))$ we have that $[i, j] \in \mu_{\text{SC}}(r_x(P))$.

We are now going to prove that $r_x(\text{next}(P)) \supseteq \mu_{SC}(r_x(P))$ and, combining with the previous result, that $r_x(\text{next}(P)) = \mu_{SC}(r_x(P))$. Consider $[i, j] \in \mu_{SC}(r_x(P))$. This means that there exists $[i, k]$ and $[h, j]$ in $r_x(P)$ such that $[i, k] \vee [h, j] = [i, j]$. Note that this means that $i \leq h$, $k \leq j$ and $h \leq k + 1$. Since $[i, k]$ and $[h, j]$ are in $r_x(P)$ there exists two individuals $z, v \in P$ such that $z_{[i,k]} = x_{[i,k]}$ and $v_{[h,j]} = x_{[h,j]}$. Now consider the individual $y = z_{[i,k]}v_{[k+1,j]}$. Since it is obtained by the crossover z and v it is inside $\text{next}(P)$. But $r_x(\{y\})$ contains $[i, j]$ since $y_{[i,j]} = z_{[i,k]}v_{[k+1,j]} = x_{[i,k]}x_{[k+1,j]} = x_{[i,j]}$. Hence the $r_x(\text{next}(P)) \supseteq \mu_{SC}(r_x(P))$. \square

As a corollary we immediately have that:

Corollary 5.4. For all $P \in \mathbf{P}$ and for all $x \in \Sigma^n$, $\{x\} \in S_1(P)$ iff $[1, n] \in \mu_{SC}(r_x(P))$.

Consequently, we have another way of computing a property of the Čech closure defined over populations by simply using the function μ_{SC} defined on SC_n :

$$\min\{i \in \mathbb{N} \mid \{x\} \subseteq S_i(P)\} = \min\{i \in \mathbb{N} \mid [1, n] \in \mu_{SC}^i(r_x(P))\}.$$

Also, note that since the function μ_{SC} maps lower sets to lower sets, finding the first $i \in \mathbb{N}$ such that iterating μ_{SC} for i times gives a set containing $[1, n]$ is equivalent to calculating the number of iterations necessary to obtain SC_n as a result. Furthermore, it is immediate that μ_{SC} is monotone, hence we can always iterate μ_{SC} until a fixed point is reached (the monotonicity of μ_{SC} and the finiteness of SC_n assures us the this is always the case). If the fixed point reached is SC_n , then the number of steps used is the minimum number i such that $\{x\} \subseteq S_i(P)$, otherwise no such i exists.

5.1. An analysis of computational complexity

The algorithm to compute the distance between two populations P_1 and P_2 is composed of two parts that must be carried on for all $x \in P_1$ (symmetrically, also for all $x \in P_2$):

1. Computing $r_x(P_2)$.
2. Finding the fixed point of μ_{SC} .

First of all, it is necessary to note that SC_n has size $O(n^2)$. Hence, the first step can be carried on in time $O(|P_2|n^3)$ steps (for any element we assume that we are comparing two strings of length n). In the second step we have that computing the fixed point of μ_{SC} can be carried on in time $O(n^7)$ steps. This time complexity is obtained in the following way: since μ_{SC} is monotone we can have at most $|SC_n|$ steps before reaching a fixed point. Every set U we are managing is of size at most $|SC_n|$ and the computation of μ_{SC} considers a comparison (in time $O(n)$) of all the pairs of U . This means that we are doing at most $O(n^2)$ iterations, where every iteration consists of $O(n^4)$ operations that can be performed in time $O(n)$. Since these operations must be carried on for all elements of P_1 , assuming all populations of size bounded by a certain constant m we have that the total computational time is $O(m^2n^3 + mn^7)$.

The time complexity bounds can be made more strict by considering two facts. The first one is that every lower set is defined by its set of maximal elements (that forms an antichain, i.e., a set such that every pair of elements is not comparable). The size of the maximal antichain of SC_n is n (its elements are $[1, 1]$, $[2, 2]$, \dots , $[n, n]$). Therefore, we can consider only sets of size $O(n)$ instead of size $O(n^2)$ and reduce the number of comparisons to $O(n^2)$ instead of $O(n^4)$. Furthermore, we may notice that if the previously cited maximal antichain is not inside $r_x(P)$, we cannot obtain $[1, n]$ into the fixed point, hence we can consider only computations when the representation of a population contains the maximal antichain. Note that when we have a set with the maximal antichain, the first iteration will contain all the elements in the form $[i, i + 1]$, the second iteration all the sets in the form $[i, i + 2]$ and, more generally, the j^{th} iteration will contains all the sets in the form $[i, i + 2^{j-1}]$. Since μ_{SC} applied to a lower set gives a lower set, we have that a fixed point is reached in $O(\log(n))$ iterations (instead of our previous bound of $O(n^2)$). Thus, the total running time can then be bounded by $O(m^2n^3 + mn^3 \log(n))$.

6. Final remarks

In this paper a crossover-based distance for genetic algorithms (GAs) has been defined. Furthermore, an algorithm of polynomial complexity in the population size and individuals length to compute this distance has been introduced. The novelty of the proposed approach consists of the following points:

- the defined distance is between populations (from which it is straightforward to obtain a family of distances between individuals), which makes modeling crossover easier;
- the representation of the GA dynamics by means of iteration of a Čech closure;
- the mathematical tools used for representing populations in our model (lower sets of a lattice).

The proposed distance could be applied to many different scenarios in GA. For example it can be of help in determining problem difficulty when used for computing the fitness distance correlation. Also, it can improve the performances of GA when used as the distance for fitness sharing. There are also other applications specific to distances between populations. For example we can try to quantify the “quality” of the genetic material of a population by computing its distance to a set of other populations. A low average distance means that the genetic material in the population is “good” (i.e., it is easy to generate new individuals).

Future work is focused on the extension of this distance to other kinds of crossover, also with the long term goal of extending it to a wider range of evolutionary algorithms (EAs). In particular, a general (representation-independent) way of extending and computing this distance should be devised, in order to provide a coherent framework for the analysis of the EAs dynamics.

Acknowledgement

This research has been supported by University of Milano Bicocca – FAR Project 2010 – ATE-0109 “Natural Computing Models with Applications to Systems Biology”.

References

- [1] J.H. Holland, *Adaptation in Natural and Artificial Systems*, The University of Michigan Press, Ann Arbor, Michigan, 1975.
- [2] D.E. Goldberg, *Genetic Algorithms in Search, Optimization and Machine Learning*, Addison-Wesley, 1989.
- [3] S. Gustafson, A. Ekárt, E. Burke, G. Kendall, Problem difficulty and code growth in genetic programming, *Genetic Programming and Evolvable Hardware* 5 (3) (2004) 271–290.
- [4] E. Burke, S. Gustafson, G. Kendall, Diversity in genetic programming: an analysis of measures and correlation with fitness, *IEEE Transactions on Evolutionary Computation* 8 (1) (2004) 47–62.
- [5] M. Tomassini, L. Vanneschi, F. Fernández, G. Galeano, A study of diversity in multipopulation genetic programming, in: 6th International Conference on Evolutionary Computation EA'03, 2003, pp. 69–81.
- [6] A. Ekárt, S.Z. Németh, Maintaining the diversity of genetic programs, in: J.A. Foster, E. Lutton, J. Miller, C. Ryan, A.G.B. Tettamanzi (Eds.), *Genetic Programming, Proceedings of the 5th European Conference, EuroGP 2002*, in: LNCS, vol. 2278, Springer, Berlin, Heidelberg, New York, Kinsale, Ireland, 2002, pp. 162–171.
- [7] T. Jones, S. Forrest, Fitness distance correlation as a measure of problem difficulty for genetic algorithms, in: L. Eshelman (Ed.), *Proceedings of the Sixth International Conference on Genetic Algorithms*, Morgan Kaufmann, San Francisco, CA, 1995, pp. 184–192.
- [8] M. Tomassini, L. Vanneschi, P. Collard, M. Clergue, A study of fitness distance correlation as a difficulty measure in genetic programming, *Evolutionary Computation* 13 (2) (2005) 213–239.
- [9] L. Vanneschi, *Theory and practice for efficient genetic programming*, Ph.D. Thesis, Faculty of Sciences, University of Lausanne, Switzerland (2004).
- [10] J. McDermott, U. O'Reilly, L. Vanneschi, K. Veeramachaneni, How far is it from here to there? A distance that is coherent with GP operators, in: S. Silva, J.A. Foster, M. Nicolau, M. Giacobini, P. Machado (Eds.), *Proceedings of the 14th European Conference on Genetic Programming, EuroGP 2011*, in: LNCS, vol. 6621, Springer Verlag, Turin, Italy, 2011, pp. 191–202.
- [11] M. Vose, *The Simple Genetic Algorithm: Foundations and Theory*, MIT Press, Cambridge, MA, USA, 1998.
- [12] M. Vose, Course notes: genetic algorithm theory, in: M. Pelikan, J. Branke (Eds.), *GECCO (Companion)*, ACM, 2010, pp. 2647–2660.
- [13] J. Rowe, M. Vose, A. Wright, Representation invariant genetic operators, *Evolutionary Computation* 18 (4) (2010) 635–660.
- [14] J. Rowe, M. Vose, A. Wright, Neighborhood graphs and symmetric genetic operators, in: *FOGA*, 2007, pp. 110–122.
- [15] H.-G. Beyer, T. Jansen, C. Reeves, M. Vose (Eds.), *Theory of Evolutionary Algorithms*, 15.–20. February 2004, in: *Dagstuhl Seminar Proceedings, Internationales Begegnungs- und Forschungszentrum für Informatik (IBFI)*, vol. 04081, Schloss Dagstuhl, Germany, 2006.
- [16] D. Arnold, T. Jansen, M.D. Vose, J. Rowe (Eds.), *Theory of Evolutionary Algorithms*, 05.02. – 10.02.2006, in: *Dagstuhl Seminar Proceedings, Internationales Begegnungs- und Forschungszentrum fuer Informatik (IBFI)*, vol. 06061, Schloss Dagstuhl, Germany, 2006.
- [17] D. Arnold, T. Jansen, J. Rowe, M. Vose, 06061 executive summary – theory of evolutionary algorithms, in: D. Arnold, T. Jansen, M. Vose, J. Rowe (Eds.), *Theory of Evolutionary Algorithms*, in: *Dagstuhl Seminar Proceedings, Internationales Begegnungs- und Forschungszentrum fuer Informatik (IBFI)*, vol. 06061, Schloss Dagstuhl, Germany, 2006.
- [18] C. Reeves, J. Rowe, *Genetic Algorithms: Principles and Perspectives: A Guide to GA Theory*, Springer, 2002.
- [19] B. Stadler, P. Stadler, M. Shpak, G. Wagner, Recombination spaces, metrics, and pretopologies, *Zeitschrift für Physikalische Chemie* 216 (2002) 2002.
- [20] P. Stadler, G. Wagner, Algebraic theory of recombination spaces, *Evolutionary Computation* 5 (3) (1997) 241–275. doi:10.1162/evco.1997.5.3.241.
- [21] G. Wagner, P. Stadler, Complex adaptations and the structure of recombination spaces, in: *School of Mathematics, UEA, Norwich NR4 7TJ*, 1997.
- [22] A. Moraglio, R. Poli, Topological interpretation of crossover, in: *Proceedings of the Genetic and Evolutionary Computation Conference*, Springer, 2004, pp. 1377–1388.
- [23] A. Moraglio, One-point geometric crossover, in: R. Schaefer, C. Cotta, J. Kolodziej, G. Rudolph (Eds.), *PPSN (1)*, in: *Lecture Notes in Computer Science*, vol. 6238, Springer, 2010, pp. 83–93.
- [24] A. Moraglio, Geometry of evolutionary algorithms, in: N. Krasnogor, P. Lanzi (Eds.), in: *GECCO (Companion)*, ACM, 2011, pp. 1439–1468.
- [25] G. Birkhoff, *Lattice theory*, American Mathematical Society, 1967.
- [26] E. Čech, *Topological Spaces*, Wiley Interscience Publisher, London, 1966.
- [27] J. Munkers, *Topology*, 2nd Edition, Prentice Hall, 1999.